

# Python 代码整理汇总-3

## Pandas\_1Foundations

DataFrame OUTLINE:  
Import data To DataFrame From Dictionary  
Import data To DataFrame From CSV file  
other import keywords: Header, Names, na\_value  
Import Big Data File  
Write into files  
DataFrame Operations: Index  
DataFrame Operations: Slicing  
DataFrame Operations: Head, Tail, Info  
DataFrame -> series -> ndarray  
DataFrame: Broadcasting  
Series and DataFrame  
Plot data in DataFrame

## 2\_ExploratoryDataAnalysis

Visual Exploratory Data Analysis  
Line Plot  
Scatter Plot  
Box Plot  
Histogram  
different plot methods  
Statistical Exploratory Data Analysis  
Summary Statistics  
Counts  
Averages  
Standard Deviations(std)  
Medians  
Quantiles  
Max, Min  
Unique

## 3\_TimeSeriesPandas

Time Serie OUTLINE:  
datetime format  
Import csv data with datetime  
slicing times  
Convert string to datetime  
Reindexing DataFrame, Filling missing values  
Resampling  
complex case  
Manipulating time series data  
String methods  
Datetime methods  
Visualization

## Pandas\_4ClimateCase

Import data file  
Add Columns  
Set Datetime as Index  
Clean numeric columns  
Data Analysis  
Statistical Features of Time Series : mean, median, max, min, sum, count  
Resample  
Practise:  
Compare Sunny and Overcast  
Visualization

## Pandas\_1Foundations

---

### DataFrame OUTLINE:

---

- Import data to DataFrame from dict & csv file
- Write into file
- DataFrame Operations: index, slicing, head, tail, info • Broadcasting
- Plot in DataFrame

### Import data To DataFrame From Dictionary

---

```

1  ## Create a DataFrame from a dictionary
2  import pandas as pd
3
4  # Pre-defined lists: country [names], Drives_Right [dr], Cars per capital[cpc]
5  names = ['United States', 'Australia', 'Japan', 'India', 'Russia', 'Morocco', 'Egypt']
6  dr = [True, False, False, False, True, True, True]
7  cpc = [809, 731, 588, 18, 200, 70, 45]
8
9  my_dict = {'country':names, 'drives_right':dr, 'cars_per_cap':cpc}
10 cars = pd.DataFrame(my_dict) # 用字典建立df
11
12 print(cars)

```

	country	drives_right	cars_per_cap
0	United States	True	809
1	Australia	False	731
2	Japan	False	588
3	India	False	18
4	Russia	True	200
5	Morocco	True	70
6	Egypt	True	45

```

1  ## Optional: use zip
2  my_dict_zip = dict(zip(['country', 'drives_right', 'cars_per_cap'], [names, dr, cpc]))
3  cars_zip = pd.DataFrame(my_dict_zip)
4  print(cars_zip)

```

	country	drives_right	cars_per_cap
0	United States	True	809
1	Australia	False	731
2	Japan	False	588
3	India	False	18
4	Russia	True	200
5	Morocco	True	70
6	Egypt	True	45

```

1  ## Specify the row labels 自定义label (列也可以)
2  row_labels = ['US', 'AUS', 'JAP', 'IN', 'RU', 'MOR', 'EG']
3  cars.index = row_labels
4  print(cars)

```

	country	drives_right	cars_per_cap
US	United States	True	809
AUS	Australia	False	731
JAP	Japan	False	588
IN	India	False	18
RU	Russia	True	200
MOR	Morocco	True	70
EG	Egypt	True	45

## Import data To DataFrame From CSV file

```

1  ## Create DataFrame from a CSV(Comma-Seperated Values) file
2  cars = pd.read_csv('cars.csv') # 分隔符默认为","
3  print(cars)
4  cars = pd.read_csv('cars.csv', index_col = 0) # 第一列为index
5  print(cars)

```

	Unnamed: 0	cars_per_cap	country	drives_right
0	US	809	United States	True
1	AUS	731	Australia	False
2	JAP	588	Japan	False
3	IN	18	India	False
4	RU	200	Russia	True
5	MOR	70	Morocco	True
6	EG	45	Egypt	True
	cars_per_cap		country	drives_right
US	809	United States		True

11	AUS	731	Australia	False
12	JAP	588	Japan	False
13	IN	18	India	False
14	RU	200	Russia	True
15	MOR	70	Morocco	True
16	EG	45	Egypt	True

## other import keywords: Header, Names, na\_value

```

1 ## header = None
2 cars = pd.read_csv('cars.csv', header=None)
3 cars

```

	0	1	2	3
0	NaN	cars_per_cap	country	drives_right
1	US	809	United States	True
2	AUS	731	Australia	False
3	JAP	588	Japan	False
4	IN	18	India	False
5	RU	200	Russia	True
6	MOR	70	Morocco	True
7	EG	45	Egypt	True

```

1 ## header = 0
2 cars = pd.read_csv('cars.csv', header=0) #第一行为header
3 cars

```

	Unnamed: 0	cars_per_cap	country	drives_right
0	US	809	United States	True
1	AUS	731	Australia	False
2	JAP	588	Japan	False
3	IN	18	India	False
4	RU	200	Russia	True
5	MOR	70	Morocco	True
6	EG	45	Egypt	True

```

1 col_names = ['A', 'B', 'C']
2 cars = pd.read_csv('cars.csv', names = col_names) # 设定列名
3 cars

```

	A	B	C
NaN	cars_per_cap	country	drives_right
US	809	United States	True
AUS	731	Australia	False
JAP	588	Japan	False
IN	18	India	False
RU	200	Russia	True
MOR	70	Morocco	True
EG	45	Egypt	True

```

1 col_names = ['A', 'B', 'C']
2 cars = pd.read_csv('cars.csv', names = col_names, na_values = {'A':['18','45']}) # A列中18, 45换成NaN
3 cars

```

	A	B	C
NaN	cars_per_cap	country	drives_right
US	809	United States	True
AUS	731	Australia	False
JAP	588	Japan	False
IN	NaN	India	False
RU	200	Russia	True
MOR	70	Morocco	True
EG	NaN	Egypt	True

## Import Big Data File

```

1 def count_entries(csv_file, c_size, colname): # 不适合一次性读取
2     """Return a dictionary with counts of
3     occurrences as value for each key."""
4
5     # Initialize an empty dictionary: counts_dict
6     counts_dict = {}
7
8     # Iterate over the file chunk by chunk
9     for chunk in pd.read_csv(csv_file, chunksize=c_size):
10         # Iterate over the column in DataFrame
11         for entry in chunk[colname]: # 仅统计指定columns
12             if entry in counts_dict.keys():
13                 counts_dict[entry] += 1
14             else:
15                 counts_dict[entry] = 1
16         # Return counts_dict
17         return counts_dict
18
19 # Call count_entries(): result_counts
20 result_counts = count_entries('cars.csv', 10, 'drives_right')
21
22 # Print result_counts
23 print(result_counts)

```

```

1 {True: 4, False: 3}

```

## Write into files

```

1 out_csv = 'newcsvFile.csv'
2 cars.to_csv(out_csv)
3 #cars.to_csv(out_csv, sep='\t')
4 #cars.to_excel('newexcelFile.xlsx')

```

## DataFrame Operations: Index

```

1 cars = pd.read_csv('cars.csv', index_col = 0)
2 print(cars)
3
4 ## show the type of dataframe
5 print(type(cars))
6 ## show the shape
7 print(cars.shape)
8 ## show the columns
9 print(cars.columns)
10 for item in cars.columns:
11     print(item)
12 ## show the index

```

```

13 print(cars.index)
14 for item in cars.index:
15     print(item)

```

```

1      cars_per_cap      country  drives_right
2  US              809  United States      True
3  AUS              731   Australia      False
4  JAP              588     Japan      False
5  IN               18      India      False
6  RU              200     Russia      True
7  MOR              70    Morocco      True
8  EG              45      Egypt      True
9  <class 'pandas.core.frame.DataFrame'>
10 (7, 3)
11 Index(['cars_per_cap', 'country', 'drives_right'], dtype='object')
12 cars_per_cap
13 country
14 drives_right
15 Index(['US', 'AUS', 'JAP', 'IN', 'RU', 'MOR', 'EG'], dtype='object')
16 US
17 AUS
18 JAP
19 IN
20 RU
21 MOR
22 EG

```

```

1  ## Show the column labels
2  print("Column Labels: " + ", ".join(cars)) # 将序列中的元素以指定的字符连接生成一个新的字符串。
3  #for col in cars:
4  #    print(col)
5
6  ## Show the row label and row content
7  for lab, row in cars.iterrows():
8      print("-----Label: ", lab)
9      print(row)

```

```

1  Column Labels: cars_per_cap, country, drives_right
2  -----Label:  US
3      cars_per_cap      809
4      country          United States
5      drives_right      True
6  Name: US, dtype: object
7  -----Label:  AUS
8      cars_per_cap      731
9      country          Australia
10     drives_right      False
11  Name: AUS, dtype: object
12 -----Label:  JAP
13     cars_per_cap      588
14     country          Japan
15     drives_right      False
16  Name: JAP, dtype: object
17 -----Label:  IN
18     cars_per_cap      18
19     country          India
20     drives_right      False
21  Name: IN, dtype: object
22 -----Label:  RU
23     cars_per_cap      200
24     country          Russia
25     drives_right      True
26  Name: RU, dtype: object
27 -----Label:  MOR
28     cars_per_cap      70
29     country          Morocco
30     drives_right      True
31  Name: MOR, dtype: object
32 -----Label:  EG
33     cars_per_cap      45
34     country          Egypt
35     drives_right      True
36  Name: EG, dtype: object

```

- for循环，短的循环放外面，长的放里面

- 调用dataframe时，存储是先行后列，for也要先行后列

## DataFrame Operations: Slicing

```
1 cars
```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False
<b>IN</b>	18	India	False
<b>RU</b>	200	Russia	True
<b>MOR</b>	70	Morocco	True
<b>EG</b>	45	Egypt	True

```
1 cars.iloc[:3,:]
2 #cars.iloc[-2:,:]
```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False

```
1 cars.loc['US']
```

```
1 cars_per_cap      809
2 country      United States
3 drives_right      True
4 Name: US, dtype: object
```

```
1 cars.loc[:, 'cars_per_cap']
```

```
1 US      809
2 AUS     731
3 JAP     588
4 IN       18
5 RU      200
6 MOR       70
7 EG       45
8 Name: cars_per_cap, dtype: int64
```

```
1 cars.cars_per_cap
```

```
1 US      809
2 AUS     731
3 JAP     588
4 IN       18
5 RU      200
6 MOR       70
7 EG       45
8 Name: cars_per_cap, dtype: int64
```

```

1 print(cars['cars_per_cap'])
2 print(type(cars['cars_per_cap']))
3
4 cars['cars_per_cap'].values

```

```

1 US      809
2 AUS     731
3 JAP     588
4 IN       18
5 RU      200
6 MOR      70
7 EG       45
8 Name: cars_per_cap, dtype: int64
9 <class 'pandas.core.series.Series'>

```

```

1 array([809, 731, 588, 18, 200, 70, 45])

```

## DataFrame Operations: Head, Tail, Info

```

1 cars

```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False
<b>IN</b>	18	India	False
<b>RU</b>	200	Russia	True
<b>MOR</b>	70	Morocco	True
<b>EG</b>	45	Egypt	True

```

1 cars.head(3) ## Default: 5

```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False

```

1 cars.tail() # 末尾

```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>JAP</b>	588	Japan	False
<b>IN</b>	18	India	False
<b>RU</b>	200	Russia	True
<b>MOR</b>	70	Morocco	True
<b>EG</b>	45	Egypt	True

```

1 cars.info()

```

```

1 <class 'pandas.core.frame.DataFrame'>
2 Index: 7 entries, US to EG
3 Data columns (total 3 columns):
4 #   Column      Non-Null Count  Dtype
5 ---  ---
6 0   cars_per_cap  7 non-null      int64
7 1   country       7 non-null      object
8 2   drives_right  7 non-null      bool
9 dtypes: bool(1), int64(1), object(1)
10 memory usage: 475.0+ bytes

```

```
1 cars.describe()
```

	<b>cars_per_cap</b>
<b>count</b>	7.000000
<b>mean</b>	351.571429
<b>std</b>	345.595552
<b>min</b>	18.000000
<b>25%</b>	57.500000
<b>50%</b>	200.000000
<b>75%</b>	659.500000
<b>max</b>	809.000000

## DataFrame -> series -> ndarray

```

1 cpc = cars['cars_per_cap']
2 print(type(cpc))
3 print(cpc.head(4))
4
5 cpc_value = cpc.values
6 print(type(cpc_value))
7 print(cpc_value)

```

```

1 <class 'pandas.core.series.Series'>
2 US      809
3 AUS     731
4 JAP     588
5 IN       18
6 Name: cars_per_cap, dtype: int64
7 <class 'numpy.ndarray'>
8 [809 731 588 18 200 70 45]

```

## DataFrame: Broadcasting

```
1 cars
```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False
<b>IN</b>	18	India	False
<b>RU</b>	200	Russia	True
<b>MOR</b>	70	Morocco	True
<b>EG</b>	45	Egypt	True



```

1 import numpy as np
2 cars.iloc[:3,0] = np.nan
3 cars

```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	NaN	United States	True
<b>AUS</b>	NaN	Australia	False
<b>JAP</b>	NaN	Japan	False
<b>IN</b>	18.0	India	False
<b>RU</b>	200.0	Russia	True
<b>MOR</b>	70.0	Morocco	True
<b>EG</b>	45.0	Egypt	True

```

1 cars['NewColumn'] = 0.0
2 cars

```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>	<b>NewColumn</b>
<b>US</b>	NaN	United States	True	0.0
<b>AUS</b>	NaN	Australia	False	0.0
<b>JAP</b>	NaN	Japan	False	0.0
<b>IN</b>	18.0	India	False	0.0
<b>RU</b>	200.0	Russia	True	0.0
<b>MOR</b>	70.0	Morocco	True	0.0
<b>EG</b>	45.0	Egypt	True	0.0

```

1 ## add column in DataFrame
2 for lab, row in cars.iterrows():
3     cars.loc[lab, "COUNTRY"] = row['country'].upper()
4 print(cars)
5
6 ## More efficient way: using .apply
7 cars['COUNTRY'] = cars['country'].apply(str.upper) # 速度更快
8 print(cars)
9
10 ## 处理数值类型的话，用np来处理更快

```

```

1      cars_per_cap  country  drives_right  NewColumn  COUNTRY
2  US           NaN  United States      True         0.0  UNITED STATES
3  AUS           NaN   Australia     False         0.0   AUSTRALIA
4  JAP           NaN     Japan      False         0.0     JAPAN
5  IN            18.0      India      False         0.0     INDIA
6  RU            200.0    Russia      True          0.0    RUSSIA
7  MOR            70.0   Morocco      True          0.0   MOROCCO
8  EG             45.0     Egypt      True          0.0     EGYPT
9      cars_per_cap  country  drives_right  NewColumn  COUNTRY
10  US           NaN  United States      True         0.0  UNITED STATES
11  AUS           NaN   Australia     False         0.0   AUSTRALIA
12  JAP           NaN     Japan      False         0.0     JAPAN
13  IN            18.0      India      False         0.0     INDIA
14  RU            200.0    Russia      True          0.0    RUSSIA
15  MOR            70.0   Morocco      True          0.0   MOROCCO
16  EG             45.0     Egypt      True          0.0     EGYPT

```

## Series and DataFrame

```

1  ## series
2  print(cars['cars_per_cap'])
3  print(type(cars['cars_per_cap']))
4
5  ## DataFrame
6  print(cars[['cars_per_cap']])
7  print(type(cars[['cars_per_cap']]))

```

```

1  US      NaN
2  AUS      NaN
3  JAP      NaN
4  IN       18.0
5  RU      200.0
6  MOR      70.0
7  EG       45.0
8  Name: cars_per_cap, dtype: float64
9  <class 'pandas.core.series.Series'>
10 cars_per_cap
11 US      NaN
12 AUS      NaN
13 JAP      NaN
14 IN       18.0
15 RU      200.0
16 MOR      70.0
17 EG       45.0
18 <class 'pandas.core.frame.DataFrame'>

```

## Plot data in DataFrame

```

1  import matplotlib.pyplot as plt
2  cars = pd.read_csv('cars.csv', index_col = 0)
3  cars

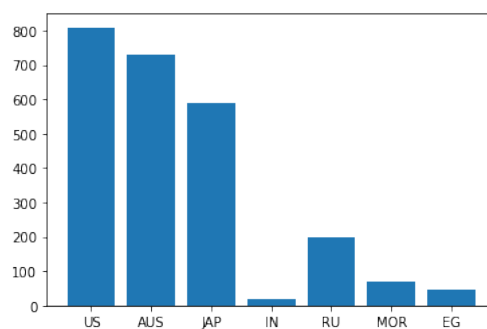
```

	cars_per_cap	country	drives_right
US	809	United States	True
AUS	731	Australia	False
JAP	588	Japan	False
IN	18	India	False
RU	200	Russia	True
MOR	70	Morocco	True
EG	45	Egypt	True

```

1  ## Plot using Matplotlib
2  cpc = cars['cars_per_cap'].values
3  xlabel = cars.index.values
4  ax = plt.bar(xlabel, cpc)
5  plt.savefig('cpc.pdf')
6  plt.show()

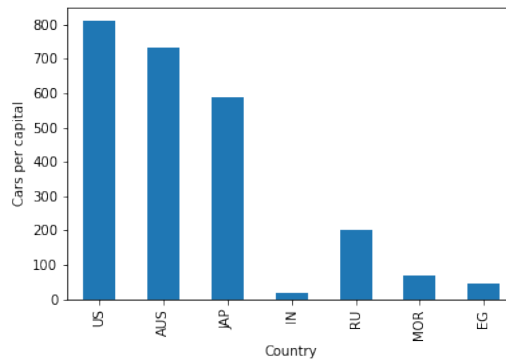
```



```

1  ## plot using pandas
2  cpc = cars['cars_per_cap']
3  ax = cpc.plot.bar(rot = 90) #标签旋转90度
4  ax.set_xlabel('Country')
5  ax.set_ylabel('Cars per capital')
6  plt.show()

```



```

1  cars

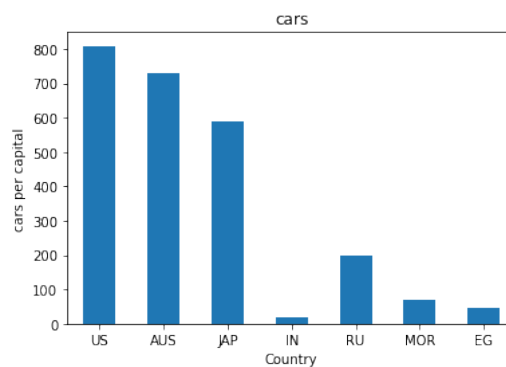
```

	<b>cars_per_cap</b>	<b>country</b>	<b>drives_right</b>
<b>US</b>	809	United States	True
<b>AUS</b>	731	Australia	False
<b>JAP</b>	588	Japan	False
<b>IN</b>	18	India	False
<b>RU</b>	200	Russia	True
<b>MOR</b>	70	Morocco	True
<b>EG</b>	45	Egypt	True

```

1  ax = cars['cars_per_cap'].plot.bar(rot=0)
2  ax.set(xlabel='Country', ylabel='cars per capital', title = 'cars')
3  plt.show()

```



## 2\_ExploratoryDataAnalysis

Example data set: **iris.csv**

iris:

- 150 observations,
- 4 features [Sepal Length, Sepal Width, Petal Length, Petal Width], • 3 species [setosa, versicolor, Virginica]

```

1  ## import data
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  iris = pd.read_csv('iris.csv')
6  col_names = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width', 'species']
7  iris.columns = col_names
8  iris.head()

```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```

1  print(iris.shape)
2  iris.info()

```

```

1  (150, 5)
2  <class 'pandas.core.frame.DataFrame'>
3  RangeIndex: 150 entries, 0 to 149
4  Data columns (total 5 columns):
5  #   Column          Non-Null Count  Dtype
6  ---  ---
7  0    sepal_length    150 non-null    float64
8  1    sepal_width     150 non-null    float64
9  2    petal_length    150 non-null    float64
10   3    petal_width     150 non-null    float64
11   4    species         150 non-null    object
12  dtypes: float64(4), object(1)
13  memory usage: 6.0+ KB

```

## Visual Exploratory Data Analysis

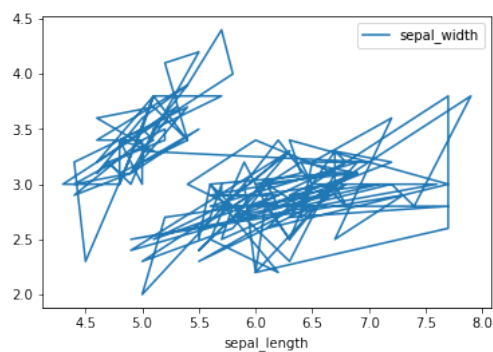
Line, Scatter, Box, Histogram

### Line Plot

```

1  iris.plot(x='sepal_length', y = 'sepal_width')
2  plt.show()

```

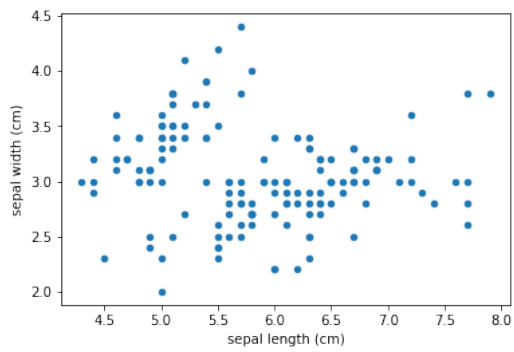


### Scatter Plot

```

1  iris.plot(x='sepal_length', y = 'sepal_width', kind='scatter')
2  plt.xlabel('sepal length (cm)')
3  plt.ylabel('sepal width (cm)')
4  plt.show()

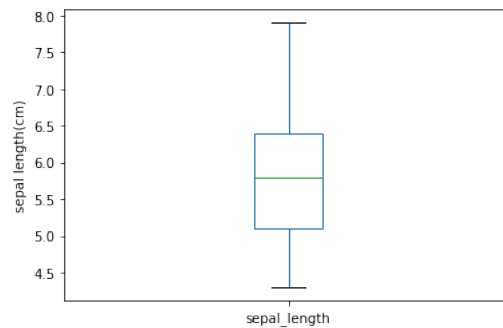
```



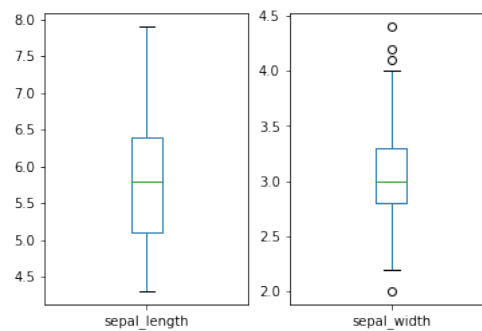
## Box Plot

[Max, third quartile 75%, median, quartile 25%, Min]

```
1 iris.plot(y = 'sepal_length', kind='box')
2 plt.ylabel('sepal length(cm)')
3 plt.show()
```

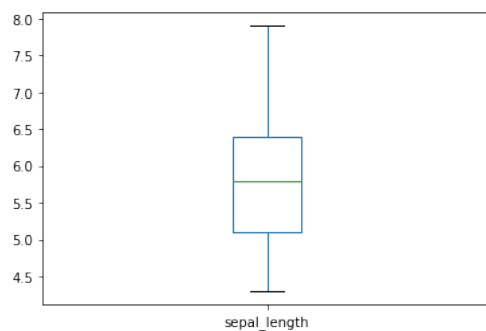


```
1 iris[['sepal_length', 'sepal_width']].plot(kind='box', subplots=True)
2 plt.show()
```



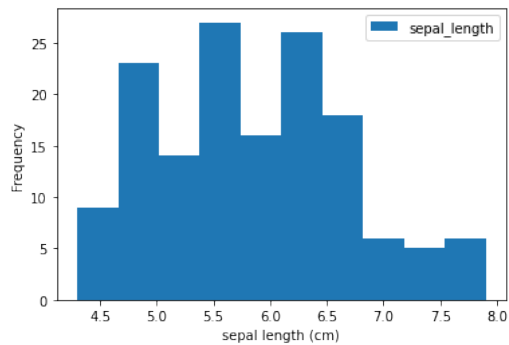
```
1 iris.sepal_length.plot(kind='box')
```

```
1 <AxesSubplot:>
```

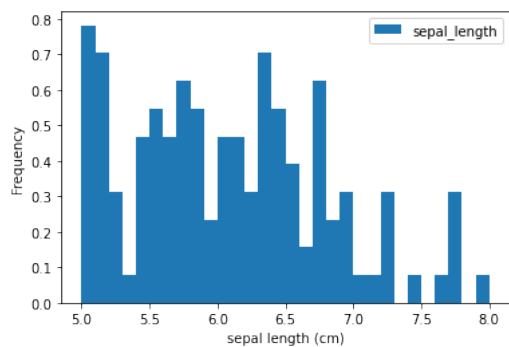


## Histogram

```
1 iris.plot(y = 'sepal_length', kind='hist')
2 plt.xlabel('sepal length (cm)')
3 plt.show()
```



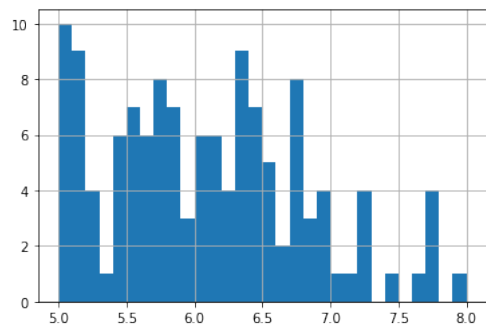
```
1 ## Options for histogram
2 # bins: integer, number of bins
3 # range: tuple, (min, max)
4 # density: boolean, normalized to 1
5 # cumulative: boolean, Compute Cumulative Distribution Function
6
7 iris.plot(y = 'sepal_length', kind='hist', bins = 30, range=(5,8), density=True)
8 plt.xlabel('sepal length (cm)')
9 plt.show()
```



## different plot methods

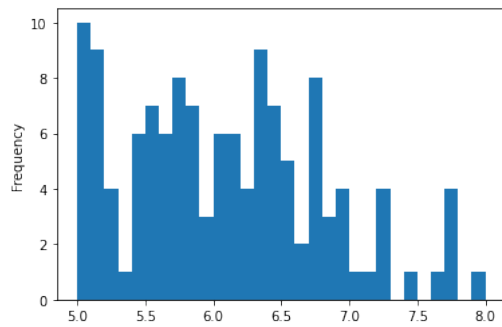
```
1 iris['sepal_length'].hist(bins = 30, range=(5,8))
```

```
1 <AxesSubplot:>
```



```
1 iris['sepal_length'].plot(bins = 30, range=(5,8), kind='hist')
```

```
1 <AxesSubplot:ylabel='Frequency'>
```



## Statistical Exploratory Data Analysis

```
1 | iris.tail(6)
```

	sepal_length	sepal_width	petal_length	petal_width	species
144	6.7	3.3	5.7	2.5	virginica
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

## Summary Statistics

```
1 | iris.describe()
```

	sepal_length	sepal_width	petal_length	petal_width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

## Counts

```
1 | iris['sepal_length'].count()
```

```
1 | 150
```

```
1 | iris[['sepal_length', 'petal_length']].count()
```

```
1 | sepal_length    150
2 | petal_length    150
3 | dtype: int64
```

## Averages

```
1 iris['sepal_length'].mean()
```

```
1 5.843333333333335
```

```
1 iris.mean()
```

```
1 sepal_length    5.843333
2 sepal_width     3.054000
3 petal_length    3.758667
4 petal_width     1.198667
5 dtype: float64
```

## Standard Deviations(std)

```
1 iris.std()
```

```
1 sepal_length    0.828066
2 sepal_width     0.433594
3 petal_length    1.764420
4 petal_width     0.763161
5 dtype: float64
```

## Medians

```
1 iris.median()
```

```
1 sepal_length    5.80
2 sepal_width     3.00
3 petal_length    4.35
4 petal_width     1.30
5 dtype: float64
```

## Quantiles

```
1 iris.quantile(0.5)
```

```
1 sepal_length    5.80
2 sepal_width     3.00
3 petal_length    4.35
4 petal_width     1.30
5 Name: 0.5, dtype: float64
```

## Max, Min

```
1 print(iris['sepal_length'].max())
2 iris.min()
```

```
1 7.9
```

```
1 sepal_length    4.3
2 sepal_width     2.0
3 petal_length    1.0
4 petal_width     0.1
5 species         setosa
6 dtype: object
```



## Unique

```
1 iris['species'].unique()
```

```
1 array(['setosa', 'versicolor', 'virginica'], dtype=object)
```

## 3\_TimeSeriesPandas

### Time Serise OUTLINE:

- import csv data with datetime • Resampling time series data
- Upsampling
- Visulization

## datetime format

### Import csv data with datetime

```
1 import pandas as pd
2
3 sales = pd.read_csv('sales.csv')
```

```
1 sales.head()
```

	Date	Company	Product	Units
0	2015-02-26 08:57:45	Streeplex	Service	4
1	2015-02-16 12:09:19	Hooli	Software	10
2	2015-02-03 14:14:18	Initech	Software	13
3	2015-02-02 08:33:01	Hooli	Software	3
4	2015-02-25 00:29:00	Initech	Service	10

```
1 sales = pd.read_csv('sales.csv', index_col = 'Date', parse_dates=True) # 转化成时间格式
2 sales.head()
```

	Company	Product	Units
Date			
2015-02-26 08:57:45	Streeplex	Service	4
2015-02-16 12:09:19	Hooli	Software	10
2015-02-03 14:14:18	Initech	Software	13
2015-02-02 08:33:01	Hooli	Software	3
2015-02-25 00:29:00	Initech	Service	10

```
1 sales.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 DatetimeIndex: 20 entries, 2015-02-26 08:57:45 to 2015-02-21 20:41:47
3 Data columns (total 3 columns):
4 #   Column   Non-Null Count  Dtype
5 ---  ---
6 0   Company  20 non-null     object
7 1   Product  20 non-null     object
8 2   Units    20 non-null     int64
9 dtypes: int64(1), object(2)
10 memory usage: 640.0+ bytes
```

## slicing times

- alternative formats:
  - sales.loc['Februray 5, 2015']
  - sales.loc['2015-Feb-5']
  - sales.loc['2015.2.5'] # 其他日期格式
- whole month:
  - sales.loc['2015-2'] # 某一月
- whole year:
  - sales.loc['2015'] # 某一年
- slicing:
  - sales.loc['2015-2-16':'2015-2-20']

```
1 print(sales.loc['2015-02-26 08:57:45'])
2 print('-'*30)
3 print(sales.loc['2015-02-26 08:57:45', 'Company'])
4 print('-'*30)
5 print(sales.iloc[0,0])
6 print('-'*30)
7 print(sales.loc['2015-2-5']) # 只匹配Date
```

```
1          Company  Product  Units
2  Date
3  2015-02-26 08:57:45  Streeplex  Service      4
4  -----
5  Date
6  2015-02-26 08:57:45  Streeplex
7  Name: Company, dtype: object
8  -----
9  Streeplex
10 -----
11          Company  Product  Units
12 Date
13 2015-02-05 01:53:06  Acme Coporation  Software      19
14 2015-02-05 22:05:03          Hooli  Service      10
```

```
1 sales.loc['2015-2-16':'2015-2-20']
```

	Company	Product	Units
Date			
2015-02-16 12:09:19	Hooli	Software	10
2015-02-19 16:02:58	Mediacore	Service	10
2015-02-19 10:59:33	Mediacore	Hardware	16

```
1 sales.loc['2015-2-4':'2015-2-5', 'Units']
```

```
1 Date
2 2015-02-05 01:53:06      19
3 2015-02-04 21:52:45      14
4 2015-02-05 22:05:03      10
5 2015-02-04 15:36:29      13
6 Name: Units, dtype: int64
```

## Convert string to datetime

```
1 dateSTR = pd.to_datetime(['2015-02-11 20:03:08', '2015-02-11 21:23', '2015-02-11 21:31'])
2 dateSTR
```

```
1 DatetimeIndex(['2015-02-11 20:03:08', '2015-02-11 21:23:00',
2               '2015-02-11 21:31:00'],
3               dtype='datetime64[ns]', freq=None)
```

## Reindexing DataFrame, Filling missing values

Ref: `pandas.DataFrame.reindex`

```
1 | sales.head()
```

	Company	Product	Units
Date			
2015-02-26 08:57:45	Streeplex	Service	4
2015-02-16 12:09:19	Hooli	Software	10
2015-02-03 14:14:18	Initech	Software	13
2015-02-02 08:33:01	Hooli	Software	3
2015-02-25 00:29:00	Initech	Service	10

pandas中的reindex方法可以为series和dataframe添加或者删除索引。

方法：`serie.reindex()`、`dataframe.reindex()`

如果新添加的索引没有对应的值，则默认为nan。如果减少索引，就相当于一个切片操作。

```
1 | sales.reindex(dateSTR) # 没有的就显示NaN
```

	Company	Product	Units
2015-02-11 20:03:08	Initech	Software	7.0
2015-02-11 21:23:00	NaN	NaN	NaN
2015-02-11 21:31:00	NaN	NaN	NaN

```
1 | sales = sales.sort_index() # 排序
2 | sales.head()
```

	Company	Product	Units
Date			
2015-02-02 08:33:01	Hooli	Software	3
2015-02-02 20:54:49	Mediacore	Hardware	9
2015-02-03 14:14:18	Initech	Software	13
2015-02-04 15:36:29	Streeplex	Software	13
2015-02-04 21:52:45	Acme Coporation	Hardware	14

```
1 | sales.reindex(dateSTR, method = 'nearest') # 参考数据最相近的就用那个去填充, 类似fillna
```

	Company	Product	Units
2015-02-11 20:03:08	Initech	Software	7
2015-02-11 21:23:00	Initech	Software	7
2015-02-11 21:31:00	Hooli	Software	4

```
1 | sales.reindex(dateSTR, method = 'backfill') #用后面的去填充
```

	Company	Product	Units
2015-02-11 20:03:08	Initech	Software	7
2015-02-11 21:23:00	Hooli	Software	4
2015-02-11 21:31:00	Hooli	Software	4

## Resampling

- Statistical methods over different time intervals
  - mean(), sum(), count()
- Resampling Frequencies
  - 'min' or 'T' for minute, 'H', 'D', 'B' for business day, 'W', 'M', 'Q' for quarter, 'A' for year

是对原样本重新处理的一个方法，是一个对常规时间序列数据重新采样和频率转换的便捷的方法。

```
1 | sales.resample('D').mean() # D: daysales.head()
```

	Company	Product	Units
Date			
2015-02-02 08:33:01	Hooli	Software	3
2015-02-02 20:54:49	Mediacore	Hardware	9
2015-02-03 14:14:18	Initech	Software	13
2015-02-04 15:36:29	Streeplex	Software	13
2015-02-04 21:52:45	Acme Coporation	Hardware	14

```
1 | sales.resample('d').mean().max()
```

```
1 | Units    14.5
2 | dtype: float64
```

```
1 | sales.resample('w').count() # week
```

	Company	Product	Units
Date			
2015-02-08	8	8	8
2015-02-15	4	4	4
2015-02-22	5	5	5
2015-03-01	3	3	3

## complex case

```
1 | sales.loc[:, 'Units'].resample('2W').sum()
```

```
1 | Date
2 | 2015-02-08    82
3 | 2015-02-22    79
4 | 2015-03-08    15
5 | Freq: 2W-SUN, Name: Units, dtype: int64
```

## Manipulating time series data

```
1 | sales = pd.read_csv('sales.csv', parse_dates = ['Date']) # 将csv中的时间字符串转换成日期格式
2 | sales.head()
```

	Date	Company	Product	Units
0	2015-02-26 08:57:45	Streeplex	Service	4
1	2015-02-16 12:09:19	Hooli	Software	10
2	2015-02-03 14:14:18	Initech	Software	13
3	2015-02-02 08:33:01	Hooli	Software	3
4	2015-02-25 00:29:00	Initech	Service	10

## String methods

```
1 sales.Company.str.upper() # 首字母大写
2 sales.Company.head()
```

```
1 0    Streeplex
2 1      Hooli
3 2    Initech
4 3      Hooli
5 4    Initech
6 Name: Company, dtype: object
```

```
1 ## Substring matching
2 sales['Product'].str.contains('ware').head()
```

```
1 0    False
2 1     True
3 2     True
4 3     True
5 4    False
6 Name: Product, dtype: bool
```

```
1 sales.Product[sales['Product'].str.contains('ware')].head()
```

```
1 1    Software
2 2    Software
3 3    Software
4 5    Software
5 7    Software
6 Name: Product, dtype: object
```

```
1 sales.Product.str.contains('ware').sum()
```

```
1 14
```

## Datetime methods

```
1 ## daytime: year:month:day:hour:minute:second
2 sales['Date'].dt.day.head() # 提取出“日”
```

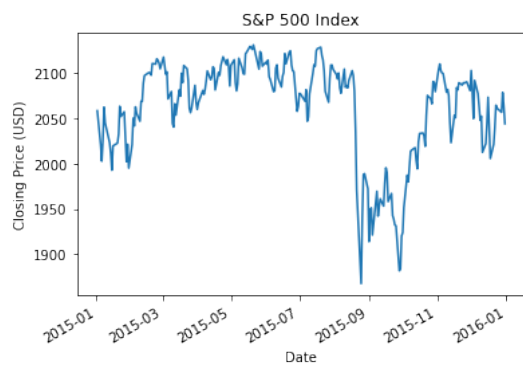
```
1 0    26
2 1    16
3 2     3
4 3     2
5 4    25
6 Name: Date, dtype: int64
```

## Visualization

```
1 import matplotlib.pyplot as plt
sp500 = pd.read_csv('sp500.csv', parse_dates = True, index_col = 0)
sp500.head()
```

	Open	High	Low	Close	Volume	Adj Close
Date						
2015-01-02	2058.899902	2072.360107	2046.040039	2058.199951	2708700000	2058.199951
2015-01-05	2054.439941	2054.439941	2017.339966	2020.579956	3799120000	2020.579956
2015-01-06	2022.150024	2030.250000	1992.439941	2002.609985	4460110000	2002.609985
2015-01-07	2005.550049	2029.609985	2005.550049	2025.900024	3805480000	2025.900024
2015-01-08	2030.609985	2064.080078	2030.609985	2062.139893	3934010000	2062.139893

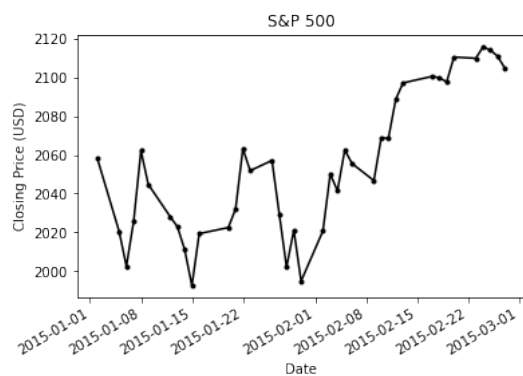
```
1 ## line plot
2 sp500['Close'].plot(title = 'S&P 500 Index')
3 plt.ylabel('Closing Price (USD)')
4 plt.show()
```



```
1 sp500.loc['2015-01-01':'2015-03-01', 'Close'].count()
```

```
1 39
```

```
1 sp500.loc['2015-01-01':'2015-03-01', 'Close'].plot(title = 'S&P 500', style = 'k.-')
2 plt.ylabel('Closing Price (USD)')
3 plt.show()
```



```
1 ## subplots
2 sp500.loc['2015', ['Close', 'Volume']].plot(subplots = True)
3 plt.show()
```



## Pandas\_4ClimateCase

### Import data file

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 df = pd.read_csv('2011_climate.txt', header=None)
6 df.head()
```

	0	1	2	3	4	5	6	7	8	9	...	34	35	36	37	38	39	40	41
0	13904	20110101	53	12	OVC045		10.00				...			29.95		AA			
1	13904	20110101	153	12	OVC049		10.00				...			30.01		AA			
2	13904	20110101	253	12	OVC060		10.00				...	030		30.01		AA			
3	13904	20110101	353	12	OVC065		10.00				...			30.03		AA			
4	13904	20110101	453	12	BKN070		10.00				...			30.04		AA			

5 rows × 44 columns

```
1 # df.info()
2 # df.columns
```

### Add Columns

```
1 ## Add columns
2 columns_label
   ='Wban,date,Time,StationType,sky_condition,sky_conditionFlag,visibility,visibilityFlag,wx_and_obst_to_vision,wx_a
   nd_obst_to_visionFlag,dry_bulb_faren,dry_bulb_farenFlag,dry_bulb_cel,dry_bulb_celFlag,wet_bulb_faren,wet_bulb_far
   enFlag,wet_bulb_cel,wet_bulb_celFlag,dew_point_faren,dew_point_farenFlag,dew_point_cel,dew_point_celFlag,relative
   _humidity,relative_humidityFlag,wind_speed,wind_speedFlag,wind_direction,wind_directionFlag,value_for_wind_charac
   ter,value_for_wind_characterFlag,station_pressure,station_pressureFlag,pressure_tendency,pressure_tendencyFlag,pr
   esschange,presschangeFlag,sea_level_pressure,sea_level_pressureFlag,record_type,hourly_precip,hourly_precipFlag,a
   ltimeter,altimeterFlag,junk'
3 columns_list = columns_label.split(',')
4 df.columns = columns_list
5 df.head()
```

	Wban	date	Time	StationType	sky_condition	sky_conditionFlag	visibility	visibilityFlag	wx_and_obst_to_vision
0	13904	20110101	53	12	OVC045		10.00		
1	13904	20110101	153	12	OVC049		10.00		
2	13904	20110101	253	12	OVC060		10.00		
3	13904	20110101	353	12	OVC065		10.00		
4	13904	20110101	453	12	BKN070		10.00		

5 rows × 44 columns

```

1  ## Remove some columns
2  columns_drop
   ='sky_conditionFlag,visibilityFlag,wx_and_obst_to_vision,wx_and_obst_to_visionFlag,dry_bulb_farenFlag,dry_bulb_celFlag,wet_bulb_farenFlag,wet_bulb_celFlag,dew_point_farenFlag,dew_point_celFlag,relative_humidityFlag,wind_speedFlag,wind_directionFlag,value_for_wind_character,value_for_wind_characterFlag,station_pressureFlag,pressure_tendencyFlag,pressure_tendencyFlag,presschange,presschangeFlag,sea_level_pressureFlag,hourly_precip,hourly_precipFlag,altimeter,record_type,altimeterFlag,junk'
3  columns_drop_list = columns_drop.split(',')
4  df_dropped = df.drop(columns_drop_list, axis = 'columns')
5  #df.head()
6  df_dropped.head()

```

	Wban	date	Time	StationType	sky_condition	visibility	dry_bulb_faren	dry_bulb_cel	wet_bulb_faren	wet_bulb_cel
0	13904	20110101	53	12	OVC045	10.00	51	10.6	38	3.1
1	13904	20110101	153	12	OVC049	10.00	51	10.6	37	3.0
2	13904	20110101	253	12	OVC060	10.00	51	10.6	37	2.9
3	13904	20110101	353	12	OVC065	10.00	50	10.0	38	3.1
4	13904	20110101	453	12	BKN070	10.00	50	10.0	37	2.8

```

1  ## OPTIONAL for columns name
2  with open('2011_climate_label.txt', 'r') as fid:
3      s = fid.read().split('\n')
4      columns_label = s[1]
5      columns_drop = s[4]
6      print(columns_label, '\n', columns_drop)

```

```

1  Wban,date,Time,StationType,sky_condition,sky_conditionFlag,visibility,visibilityFlag,wx_and_obst_to_vision,wx_and_obst_to_visionFlag,dry_bulb_faren,dry_bulb_farenFlag,dry_bulb_cel,dry_bulb_celFlag,wet_bulb_faren,wet_bulb_farenFlag,wet_bulb_cel,wet_bulb_celFlag,dew_point_faren,dew_point_farenFlag,dew_point_cel,dew_point_celFlag,relative_humidity,relative_humidityFlag,wind_speed,wind_speedFlag,wind_direction,wind_directionFlag,value_for_wind_character,value_for_wind_characterFlag,station_pressure,station_pressureFlag,pressure_tendency,pressure_tendencyFlag,presschange,presschangeFlag,sea_level_pressure,sea_level_pressureFlag,record_type,hourly_precip,hourly_precipFlag,altimeter,altimeterFlag,junk
sky_conditionFlag,visibilityFlag,wx_and_obst_to_vision,wx_and_obst_to_visionFlag,dry_bulb_farenFlag,dry_bulb_celFlag,wet_bulb_farenFlag,wet_bulb_celFlag,dew_point_farenFlag,dew_point_celFlag,relative_humidityFlag,wind_speedFlag,wind_directionFlag,value_for_wind_character,value_for_wind_characterFlag,station_pressureFlag,pressure_tendencyFlag,pressure_tendencyFlag,presschange,presschangeFlag,sea_level_pressureFlag,hourly_precip,hourly_precipFlag,altimeter,record_type,altimeterFlag,junk

```

## Set Datetime as Index

REF:

Pandas.DataFrame.astype

Pandas.to\_datetime

```

1  df_dropped.date

```

```

1  0      20110101
2  1      20110101
3  2      20110101
4  3      20110101
5  4      20110101
6  ...
7  10332   20111231
8  10333   20111231
9  10334   20111231
10 10335   20111231
11 10336   20111231
12 Name: date, Length: 10337, dtype: int64

```

```

1  ## Convert columns to pandas datetime object
2  # print(df_dropped.info())
3  # print(df_dropped.head())
4  ## convert dtypes from int64 to str(object)
5  df_dropped['date'] = df_dropped['date'].astype('str')

```



```

6
7  ## convert time -> hour+minute: 453 -> 0453 | 53 -> 0053
8  df_dropped['Time'] = df_dropped['Time'].apply(lambda x: '{:0>4}'.format(x))
9
10 ## pandas datetime object
11 date_times = pd.to_datetime(df_dropped['date'] + df_dropped['Time'], format='%Y%m%d%H%M')
12 ## set_index
13 df_clean = df_dropped.set_index(date_times)
14 df_clean.head()

```

	Wban	date	Time	StationType	sky_condition	visibility	dry_bulb_faren	dry_bulb_cel	wet_bulb_faren
<b>2011-01-01 00:53:00</b>	13904	20110101	0053	12	OVC045	10.00	51	10.6	38
<b>2011-01-01 01:53:00</b>	13904	20110101	0153	12	OVC049	10.00	51	10.6	37
<b>2011-01-01 02:53:00</b>	13904	20110101	0253	12	OVC060	10.00	51	10.6	37
<b>2011-01-01 03:53:00</b>	13904	20110101	0353	12	OVC065	10.00	50	10.0	38
<b>2011-01-01 04:53:00</b>	13904	20110101	0453	12	BKN070	10.00	50	10.0	37

```

1 df_dropped.info()

```

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 10337 entries, 0 to 10336
3 Data columns (total 17 columns):
4  #   Column                Non-Null Count  Dtype
5  ---  ---
6  0    Wban                   10337 non-null  int64
7  1    date                   10337 non-null  object
8  2    Time                   10337 non-null  object
9  3    StationType            10337 non-null  int64
10  4    sky_condition          10337 non-null  object
11  5    visibility              10325 non-null  object
12  6    dry_bulb_faren         10337 non-null  object
13  7    dry_bulb_cel           10337 non-null  object
14  8    wet_bulb_faren         10337 non-null  object
15  9    wet_bulb_cel           10337 non-null  object
16  10   dew_point_faren        10337 non-null  object
17  11   dew_point_cel          10337 non-null  object
18  12   relative_humidity      10337 non-null  object
19  13   wind_speed              10337 non-null  object
20  14   wind_direction         10337 non-null  object
21  15   station_pressure       10337 non-null  object
22  16   sea_level_pressure     10337 non-null  objectd
23  types: int64(2), object(15)memory usage: 1.3+ MB

```

## Clean numeric columns

```

1  ## Read 'dry_bulb_faren' Temperature
2  df_clean.loc['2011.6.20 8AM':'2011.6.20 9AM', 'dry_bulb_faren']

```

```

1 | 2011-06-20 08:27:00      M
2 | 2011-06-20 08:28:00      M
3 | 2011-06-20 08:29:00      M
4 | 2011-06-20 08:30:00      M
5 | 2011-06-20 08:31:00      M
6 | 2011-06-20 08:32:00      M
7 | 2011-06-20 08:33:00      M
8 | 2011-06-20 08:34:00      M
9 | 2011-06-20 08:35:00      M
10 | 2011-06-20 08:53:00      83
11 | 2011-06-20 09:08:00      84
12 | 2011-06-20 09:53:00      88
13 | Name: dry_bulb_faren, dtype: object

```

```

1 | df_clean['dry_bulb_faren'] = pd.to_numeric(df_clean['dry_bulb_faren'], errors='coerce') #coerce: set 'M' to NaN
2 | df_clean.loc['2011.6.20 8AM':'2011.6.20 9AM', 'dry_bulb_faren']

```

```

1 | 2011-06-20 08:27:00      NaN
2 | 2011-06-20 08:28:00      NaN
3 | 2011-06-20 08:29:00      NaN
4 | 2011-06-20 08:30:00      NaN
5 | 2011-06-20 08:31:00      NaN
6 | 2011-06-20 08:32:00      NaN
7 | 2011-06-20 08:33:00      NaN
8 | 2011-06-20 08:34:00      NaN
9 | 2011-06-20 08:35:00      NaN
10 | 2011-06-20 08:53:00      83.0
11 | 2011-06-20 09:08:00      84.0
12 | 2011-06-20 09:53:00      88.0
13 | Name: dry_bulb_faren, dtype: float64

```

```

1 | df_clean['wind_speed'] = pd.to_numeric(df_clean['wind_speed'], errors='coerce')df_clean['visibility'] =
  pd.to_numeric(df_clean['visibility'], errors='coerce')df_clean['dew_point_faren'] =
  pd.to_numeric(df_clean['dew_point_faren'], errors='coerce')
2 | df_clean.head()

```

	Wban	date	Time	StationType	sky_condition	visibility	dry_bulb_faren	dry_bulb_cel	wet_bulb_faren
2011-01-01 00:53:00	13904	20110101	0053	12	OVC045	10.0	51.0	10.6	38
2011-01-01 01:53:00	13904	20110101	0153	12	OVC049	10.0	51.0	10.6	37
2011-01-01 02:53:00	13904	20110101	0253	12	OVC060	10.0	51.0	10.6	37
2011-01-01 03:53:00	13904	20110101	0353	12	OVC065	10.0	50.0	10.0	38
2011-01-01 04:53:00	13904	20110101	0453	12	BKN070	10.0	50.0	10.0	37

## Data Analysis

### Statistical Features of Time Series : mean, median, max, min, sum, count

```

1 | print(df_clean['dry_bulb_faren'].median())
2 | print(df_clean.loc['2011-Apr':'2011-Jun', 'dry_bulb_faren'].max())
3 | print(df_clean.loc['2011-Jan', 'dry_bulb_faren'].count())

```

```

1 | 72.0
2 | 104.0
3 | 907

```

## Resample

```
1 daily_mean_2011 = df_clean.resample('D').mean()
2 daily_mean_2011
```

	Wban	StationType	visibility	dry_bulb_faren	dew_point_faren	wind_speed
2011-01-01	13904	12	10.000000	50.166667	20.500000	11.083333
2011-01-02	13904	12	10.000000	39.416667	19.708333	4.166667
2011-01-03	13904	12	10.000000	46.846154	35.500000	2.653846
2011-01-04	13904	12	5.071429	53.367347	50.408163	2.510204
2011-01-05	13904	12	7.672414	57.965517	40.068966	4.689655
...	...	...	...	...	...	...
2011-12-27	13904	12	10.000000	44.833333	32.125000	4.458333
2011-12-28	13904	12	10.000000	45.750000	35.166667	5.375000
2011-12-29	13904	12	10.000000	50.320000	35.600000	4.480000
2011-12-30	13904	12	10.000000	52.541667	37.875000	5.500000
2011-12-31	13904	12	10.000000	54.458333	42.416667	5.416667

365 rows × 6 columns

```
1 # dataframe -> series -> numpy array, NOTICE: values is attribute, not method
2 daily_temperature_2011 = daily_mean_2011['dry_bulb_faren'].values
3 print(type(daily_mean_2011['dry_bulb_faren']))
4 print(type(daily_temperature_2011))
```

```
1 <class 'pandas.core.series.Series'>
2 <class 'numpy.ndarray'>
```

## Practise:

- import '2010\_climate.csv' and format it
- set 'Date' column as its index
- Compare the temperature for 2011 and 2010

```
1 ## add another dataset: weather_data_2010.csv
2 df_climate = pd.read_csv('2010_climate.csv', header = 0)
3 print(df_climate.head())
4 df_climate2 = df_climate.set_index('Date')
5 print(df_climate2.head())
6 df_climate3 = df_climate2.reset_index()
7 print(df_climate3.head())
```

```
1      Temperature  DewPoint  Pressure      Date
2      0          46.2      37.5        1.0  20100101 00:00
3      1          44.6      37.1        1.0  20100101 01:00
4      2          44.1      36.9        1.0  20100101 02:00
5      3          43.8      36.9        1.0  20100101 03:00
6      4          43.5      36.8        1.0  20100101 04:00
7
8      Temperature  DewPoint  Pressure
9      Date
10 20100101 00:00      46.2      37.5        1.0
11 20100101 01:00      44.6      37.1        1.0
12 20100101 02:00      44.1      36.9        1.0
13 20100101 03:00      43.8      36.9        1.0
14 20100101 04:00      43.5      36.8        1.0
15
16      Date  Temperature  DewPoint  Pressure
17 0 20100101 00:00      46.2      37.5        1.0
18 1 20100101 01:00      44.6      37.1        1.0
19 2 20100101 02:00      44.1      36.9        1.0
20 3 20100101 03:00      43.8      36.9        1.0
21 4 20100101 04:00      43.5      36.8        1.0
```

```

1 df_climate['Date'] = pd.to_datetime(df_climate['Date'], format = '%Y%m%d %H:%M', errors='coerce')
2 df_climate = df_climate.set_index('Date')
3 df_climate.head()

```

	Temperature	DewPoint	Pressure
Date			
2010-01-01 00:00:00	46.2	37.5	1.0
2010-01-01 01:00:00	44.6	37.1	1.0
2010-01-01 02:00:00	44.1	36.9	1.0
2010-01-01 03:00:00	43.8	36.9	1.0
2010-01-01 04:00:00	43.5	36.8	1.0

```

1 daily_climate = df_climate.resample('D').mean()
2 daily_climate
3
4 # daily_temperature_climate = daily_climate.reset_index()
5 # daily_temperature_climate.head()

```

	Temperature	DewPoint	Pressure
Date			
2010-01-01	49.337500	37.716667	1.0
2010-01-02	49.795833	38.370833	1.0
2010-01-03	49.900000	38.279167	1.0
2010-01-04	49.729167	38.008333	1.0
2010-01-05	49.841667	38.087500	1.0
...	...	...	...
2010-12-27	49.204167	37.816667	1.0
2010-12-28	48.979167	37.329167	1.0
2010-12-29	48.804167	37.025000	1.0
2010-12-30	49.008333	37.325000	1.0
2010-12-31	49.195833	37.450000	1.0

365 rows × 3 columns

```

1 daily_temperature_2010 = daily_climate['Temperature']
2 print(np.mean(daily_temperature_2011 - daily_temperature_2010))

```

```

1 1.3301831870056482

```

## Compare Sunny and Overcast

```

1 is_sky_clear = df_clean['sky_condition'] == 'CLR'
2 sunny = df_clean[is_sky_clear]
3 sunny.head()

```

	Wban	date	Time	StationType	sky_condition	visibility	dry_bulb_faren	dry_bulb_cel	wet_bulb_faren
2011-01-01 13:53:00	13904	20110101	1353	12	CLR	10.0	59.0	15.0	45
2011-01-01 14:53:00	13904	20110101	1453	12	CLR	10.0	59.0	15.0	45
2011-01-01 15:53:00	13904	20110101	1553	12	CLR	10.0	57.0	13.9	44
2011-01-01 16:53:00	13904	20110101	1653	12	CLR	10.0	55.0	12.8	43
2011-01-01 17:53:00	13904	20110101	1753	12	CLR	10.0	50.0	10.0	40

```

1 sunny_daily_max = sunny.resample('D').max()
2 sunny_daily_max.head()

```

	Wban	date	Time	StationType	sky_condition	visibility	dry_bulb_faren	dry_bulb_cel	wet_bulb_faren
2011-01-01	13904.0	20110101	2353	12.0	CLR	10.0	59.0	8.3	45
2011-01-02	13904.0	20110102	2253	12.0	CLR	10.0	35.0	1.7	32
2011-01-03	13904.0	20110103	0453	12.0	CLR	10.0	32.0	0.0	29
2011-01-04	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2011-01-05	13904.0	20110105	2353	12.0	CLR	10.0	35.0	1.7	33

```

1 is_sky_overcast = df_clean['sky_condition'].str.contains('OVC')
2 overcast = df_clean.loc[is_sky_overcast]
3 overcast_daily_max = overcast.resample('D').max()

```

```

1 # Calculate the mean of sunny_daily_max
2 sunny_daily_max_mean = sunny_daily_max.mean()
3
4 # Calculate the mean of overcast_daily_max
5 overcast_daily_max_mean = overcast_daily_max.mean()
6
7 # Print the difference (sunny minus overcast)
8 print(sunny_daily_max_mean - overcast_daily_max_mean)

```

```

1 Wban          0.000000
2 StationType   0.000000
3 visibility     0.174359
4 dry_bulb_faren 6.504304
5 dew_point_faren -4.339286
6 wind_speed     -3.246062
7 dtype: float64

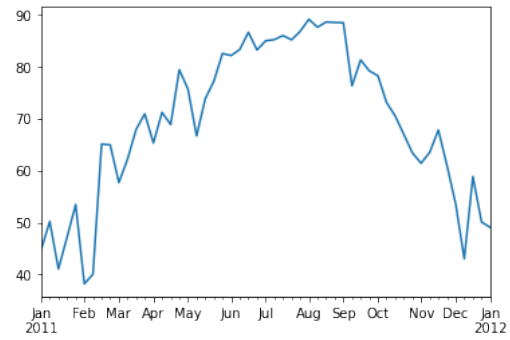
```

## Visualization

```

1 ## Show the visibility and dry_bulb_faren in figure
2 weekly_mean = df_clean['dry_bulb_faren'].resample('W').mean()
3 weekly_mean.plot()
4 plt.show()

```



```

1 monthly_max = df_clean[['dew_point_faren', 'dry_bulb_faren']].resample('M').max()
2 monthly_max.plot(kind='hist', bins=8, alpha=0.5, subplots=True)
3 plt.show()

```

